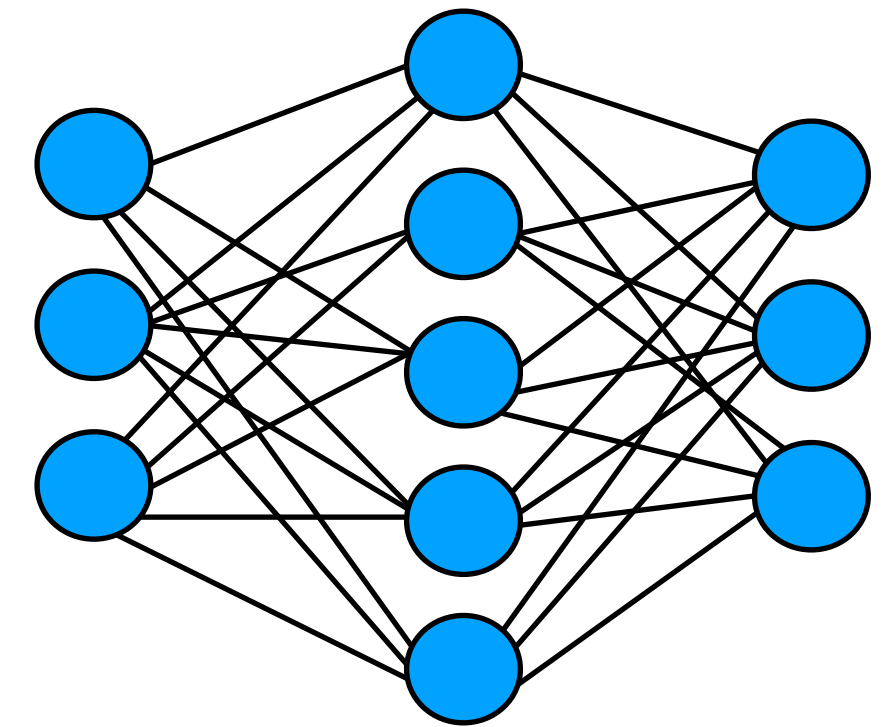
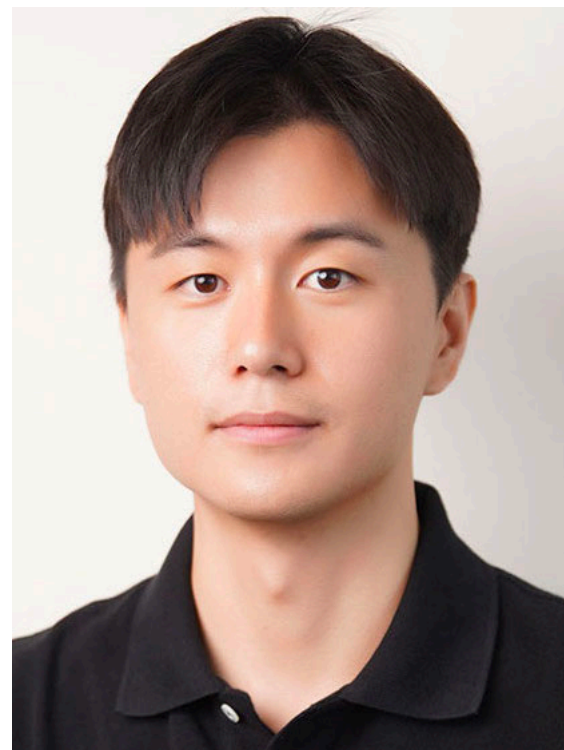


A Representer Theorem for Vector-Valued Neural Networks: Insights on Weight Decay Regularization and Widths of DNNs

Joe Shenouda
UW-Madison ECE



Rahul Parhi



Kangwook Lee



Rob Nowak

DP4ML ICML 2023

Motivations

- **Vector valued NNs** are key to understanding deep networks and multi-task learning.
- **Weight decay** is the most popular explicit regularizer for training deep neural networks (DNNs) and can have a **drastic effect** on the generalization ability of the network.

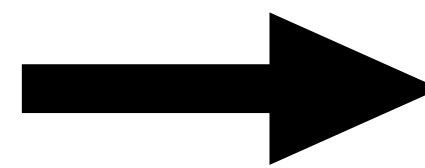
Background

$$f_{\theta} = \text{[Diagram of a fully connected neural network with 3 input nodes, 4 hidden nodes, and 3 output nodes]} \quad f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^D$$

- Training a DNN with **weight decay** corresponds to minimizing a data fidelity loss plus the sum of squared weights.

Training with Weight Decay

$$\theta^{k+1} = \theta^k - \gamma \nabla_{\theta^k} \mathcal{L} - \boxed{\gamma \lambda \theta^k}$$



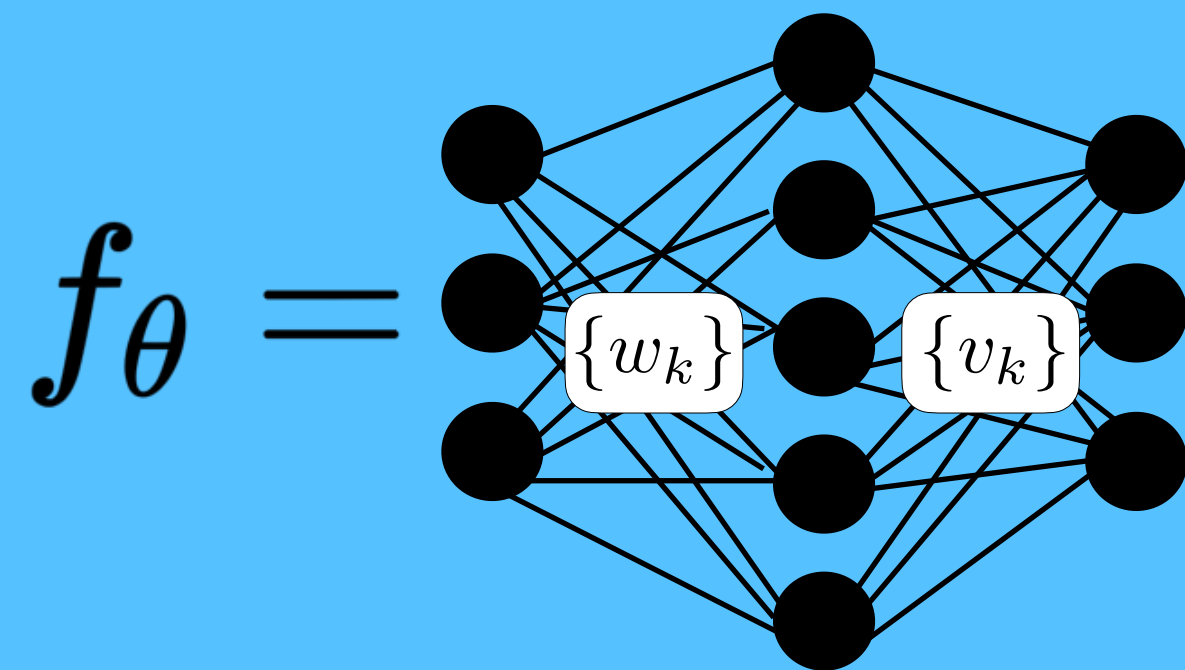
Weight Decay Objective

$$\min_{\theta} \sum_{i=1}^N \mathcal{L}(\mathbf{y}_i, f(\mathbf{x}_i)) + \frac{\lambda}{2} \|\theta\|_2^2$$

Neural Balance Theorem

Theorem (NBT [Yang 2022, Parhi 2023])

Let f_θ be a function represented by a homogenous DNN that **solves the weight decay objective**.

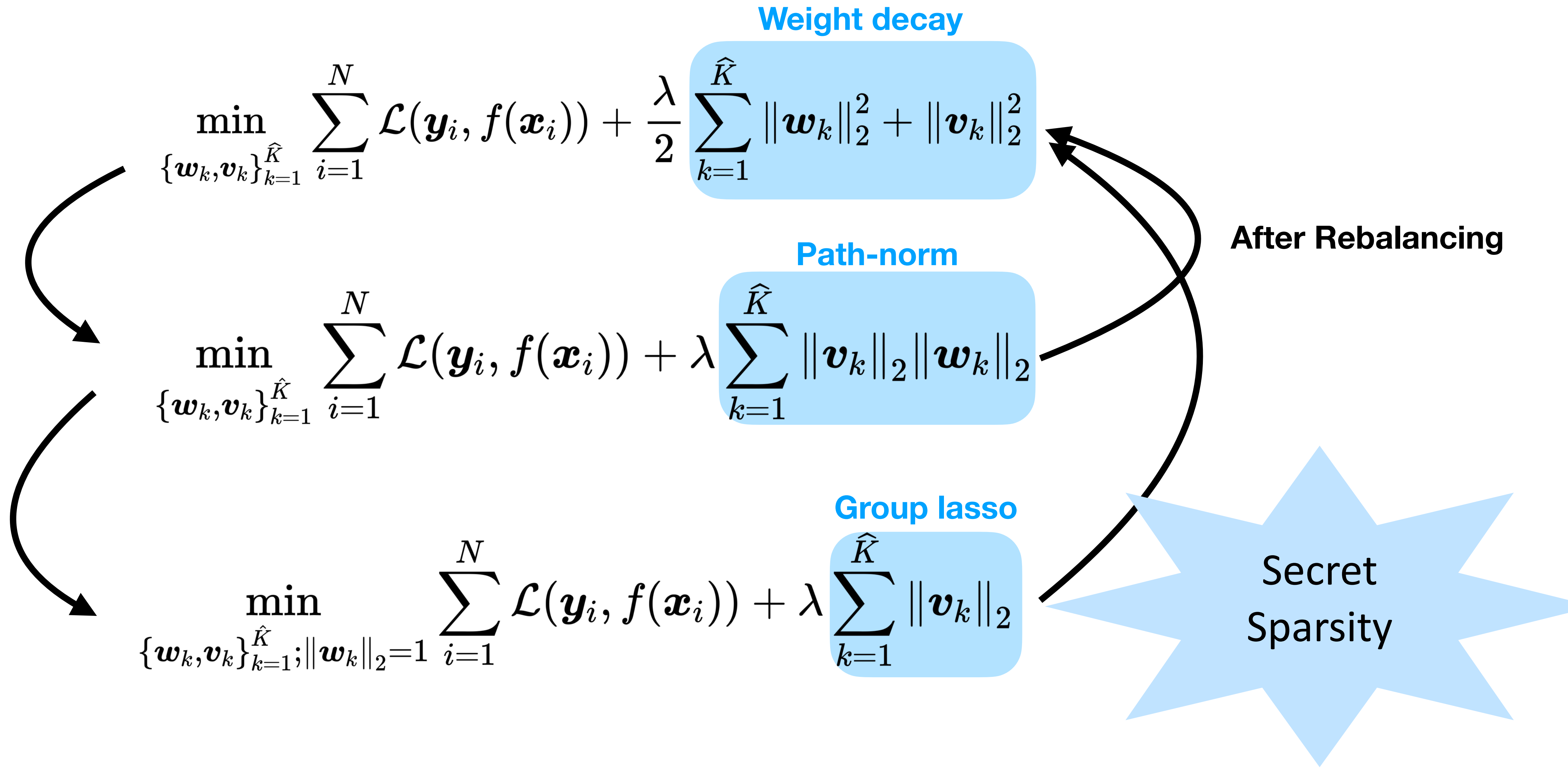


$$\min_{\{w_k, v_k\}} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i)) + \frac{\lambda}{2} \sum_{k=1}^K (\|w_k\|_2^2 + \|v_k\|_2^2)$$

Then for any **neuron** with input weight w_k and output weight v_k we have $\|w_k\|_2 = \|v_k\|_2$

Balanced!

The Secret Sparsity of Weight Decay



Outline of Contributions

Vector-Valued Variation Spaces (VV Spaces)

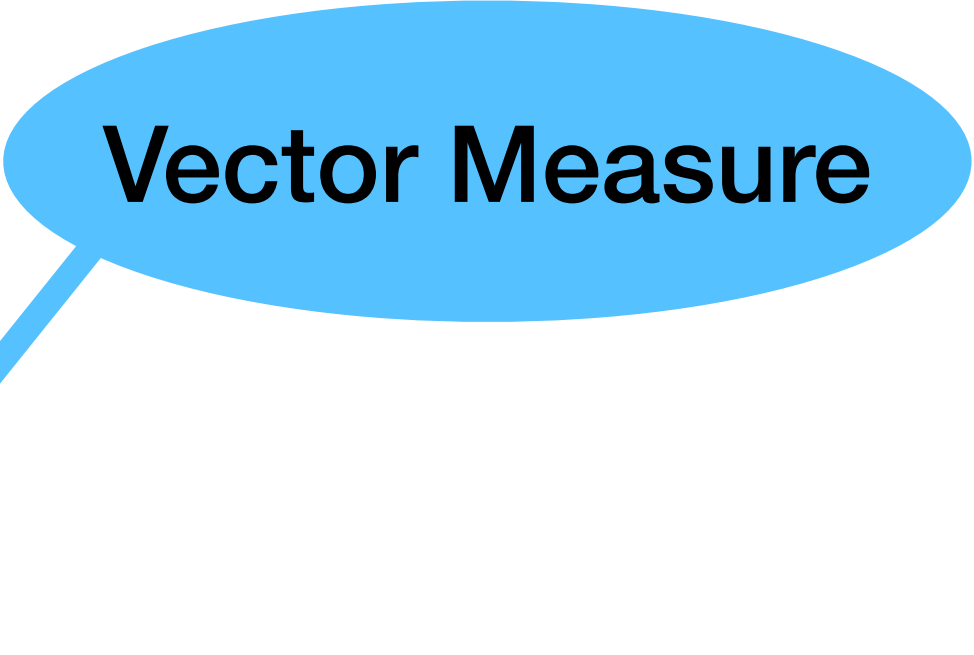
- **Representer Theorem** showing shallow vector-valued NNs trained with weight decay solve data fitting problem over the VV Space.

Bounds on Necessary Width

- **Tighter bounds** on the necessary width for any homogenous layer of DNN, depending only on the **rank** of the internal features.

W Spaces

Vector Valued Variation Spaces (VV Spaces)

$$VV = \left\{ f(\mathbf{x}) = \int_{S^d} \sigma(\mathbf{w}^T \mathbf{x}) d\nu(\mathbf{w}) : \|f\|_{VV} < \infty \right\}$$


$$\|f\|_{VV} = \sum_k \|v_k\|_2 \quad \text{(For finite-width networks)}$$

Representer Theorem

Also extends to deep neural networks

For any dataset $\{x_i, y_i\}_{i=1}^N$ and any lower semicontinuous loss function \mathcal{L} there exists a solution to,

$$\arg \min_{f \in VV} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i)) + \lambda \|f\|_{VV}$$

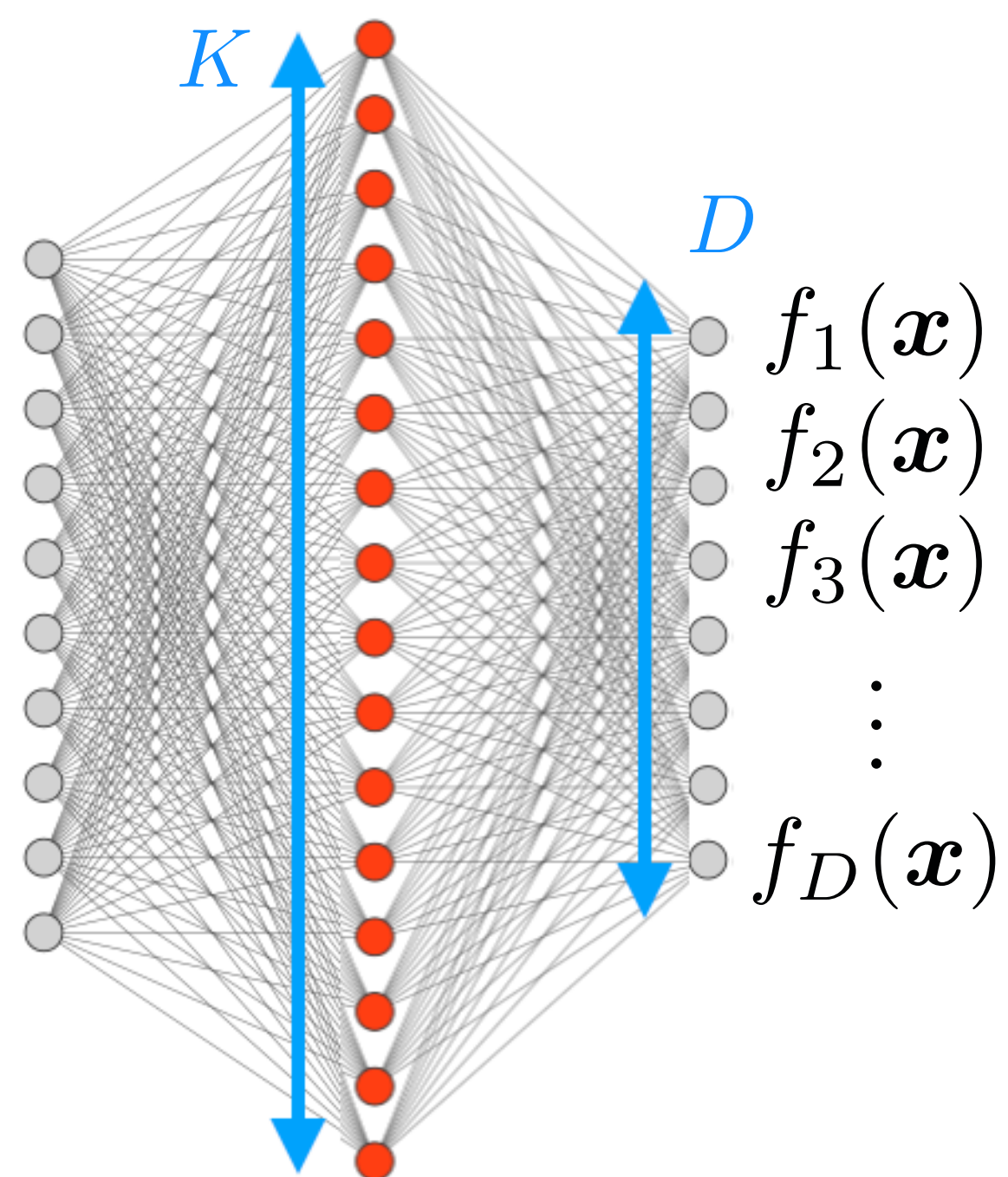
with the following **representation**

$$f_{\theta^*}(x) = \sum_{k=1}^{K_0} \mathbf{v}_k \sigma(\mathbf{w}_k^T \mathbf{x})$$

$$K_0 \leq N^2$$

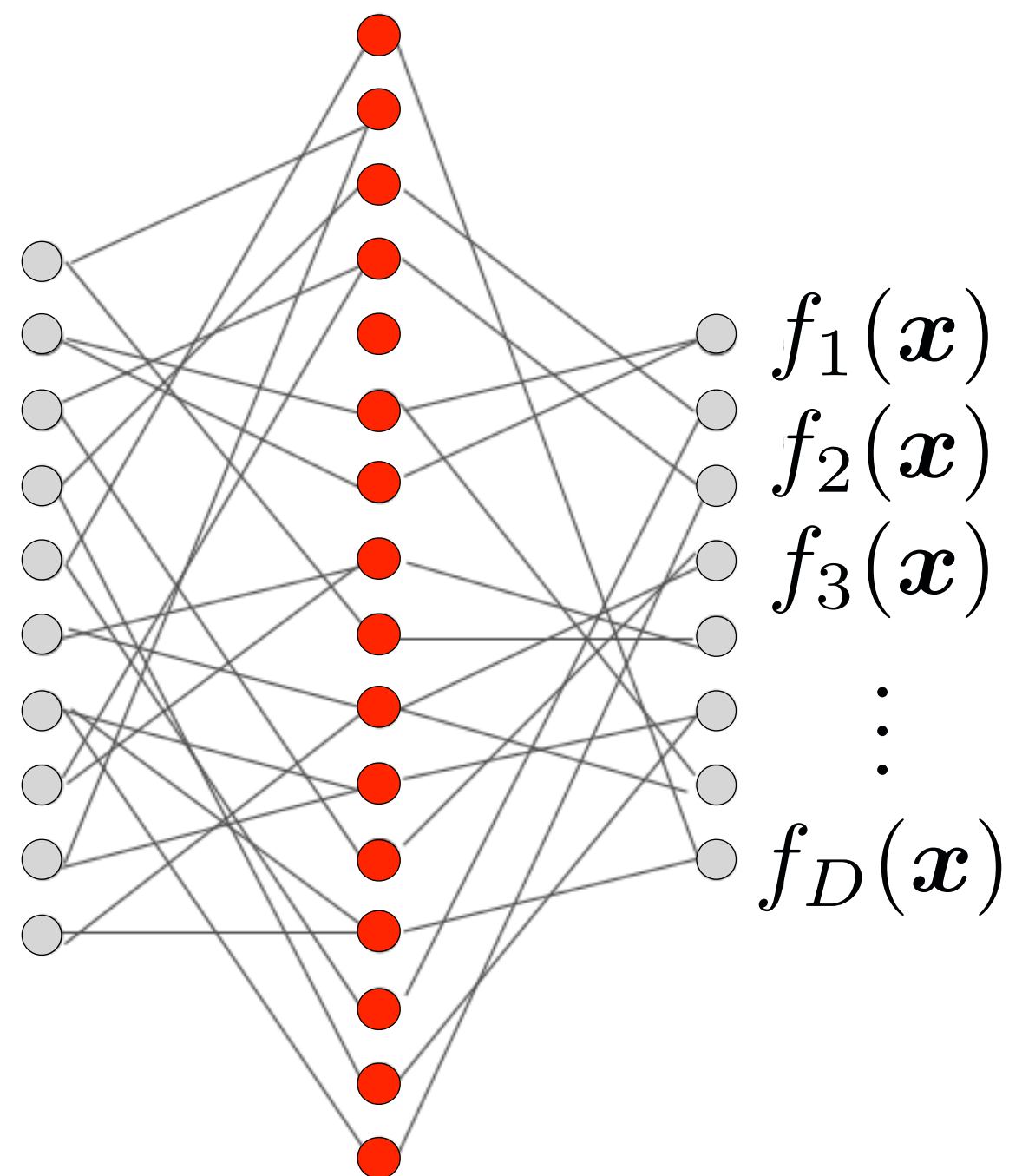
The dual space of vector-valued continuous functions is the space of vector-valued measures.

dense weights



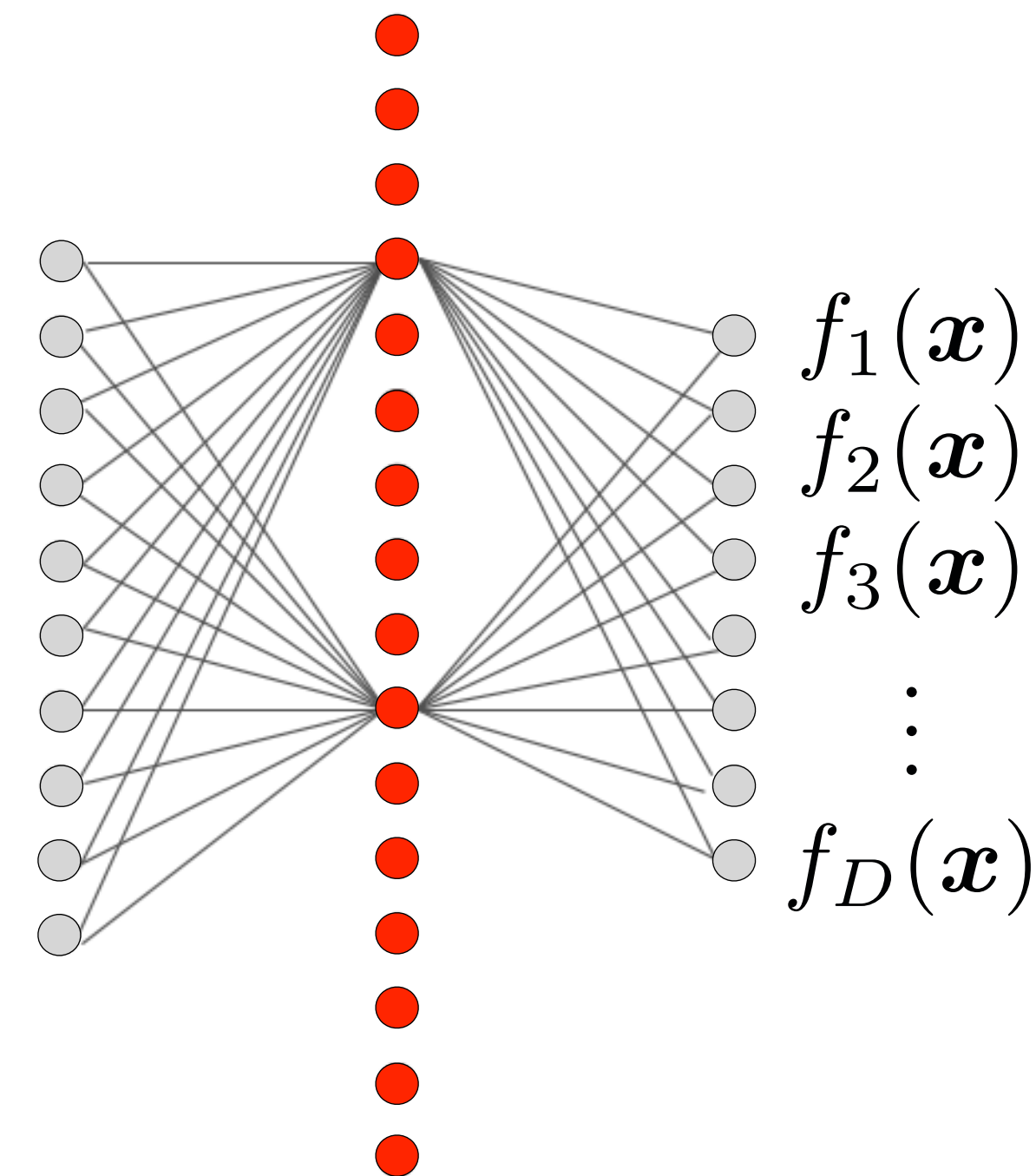
$$\|f\|_{VV} = O(K\sqrt{D})$$

sparse weights



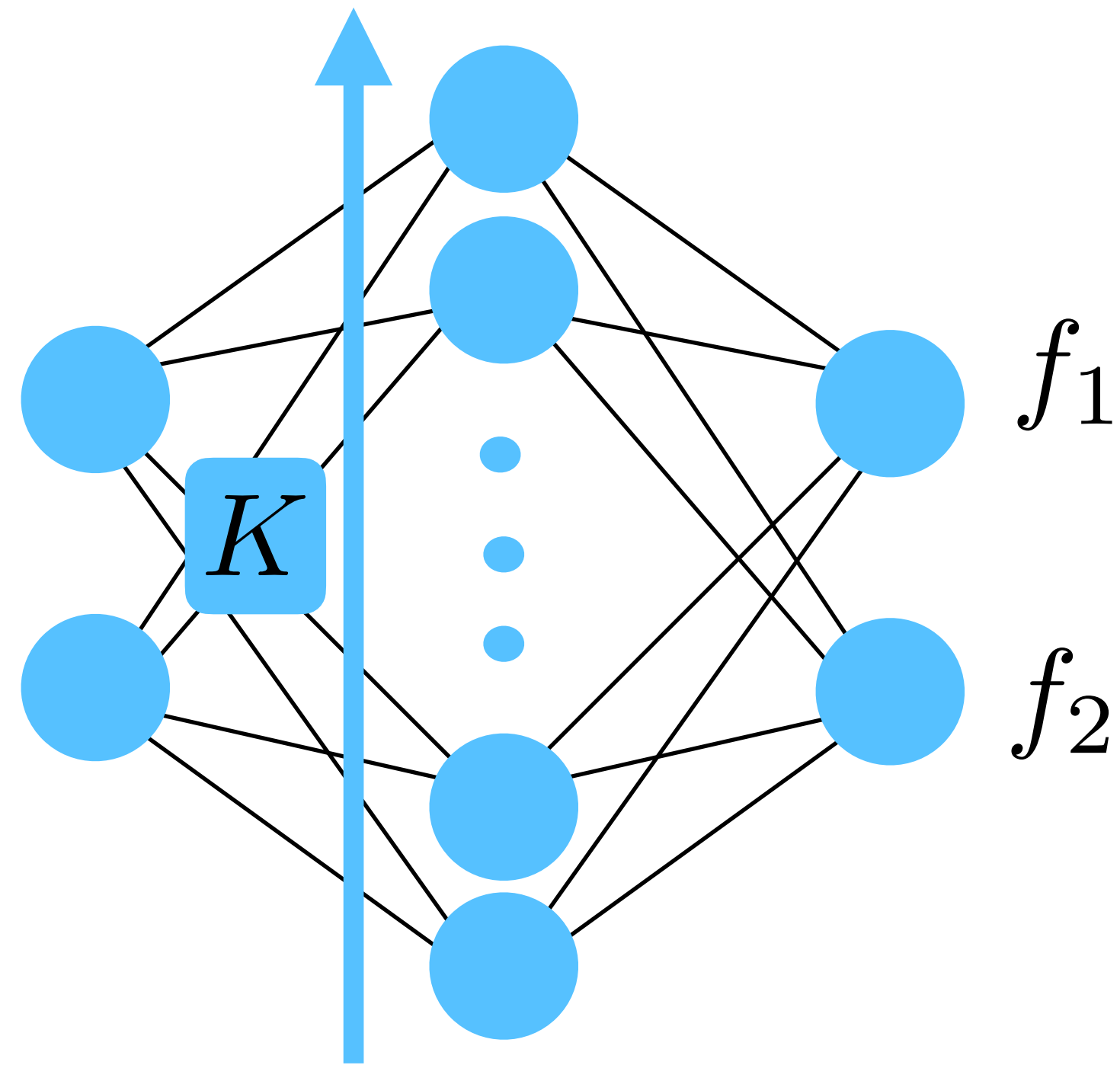
$$\|f\|_{VV} = O(D)$$

sparse neurons



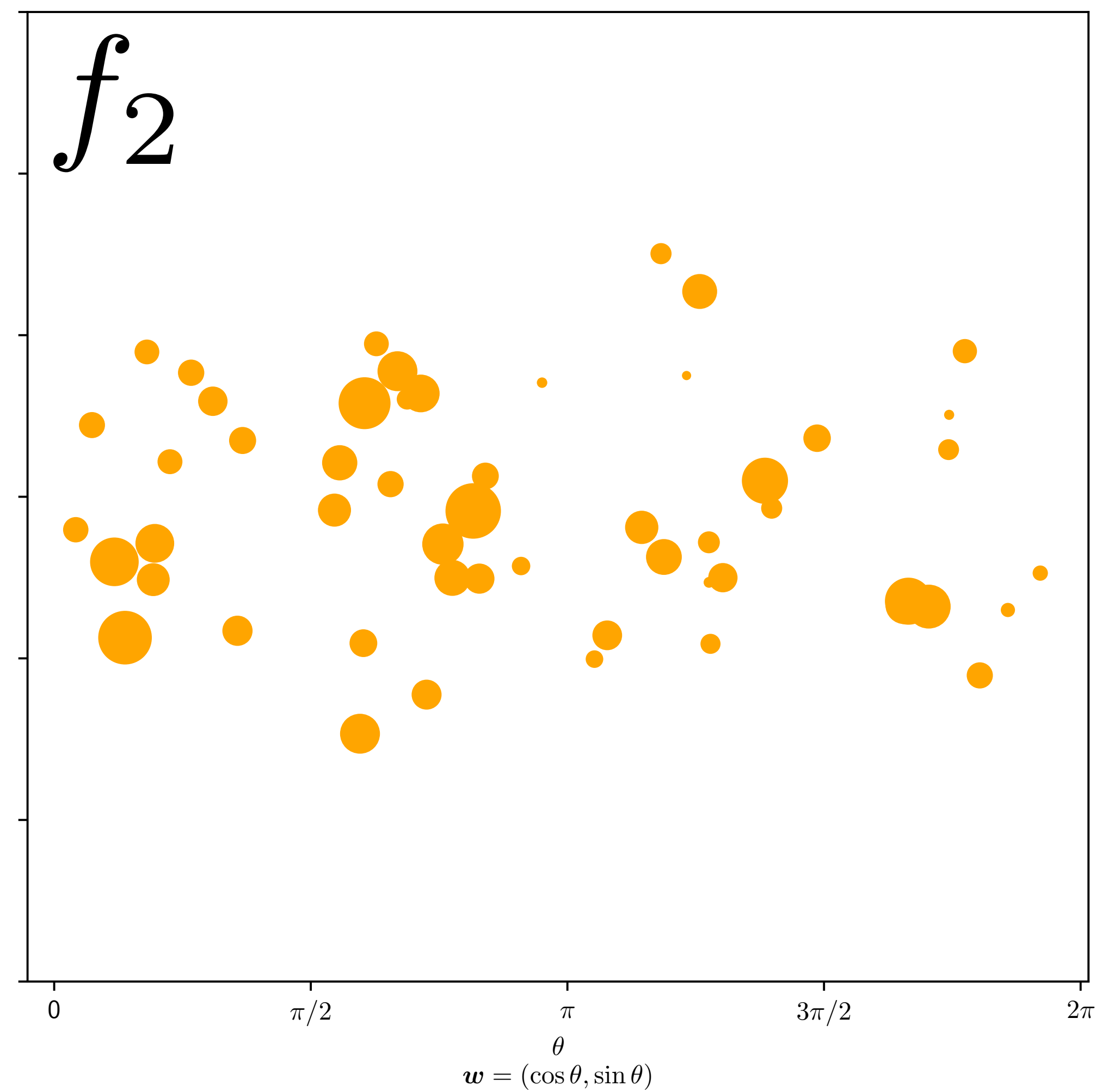
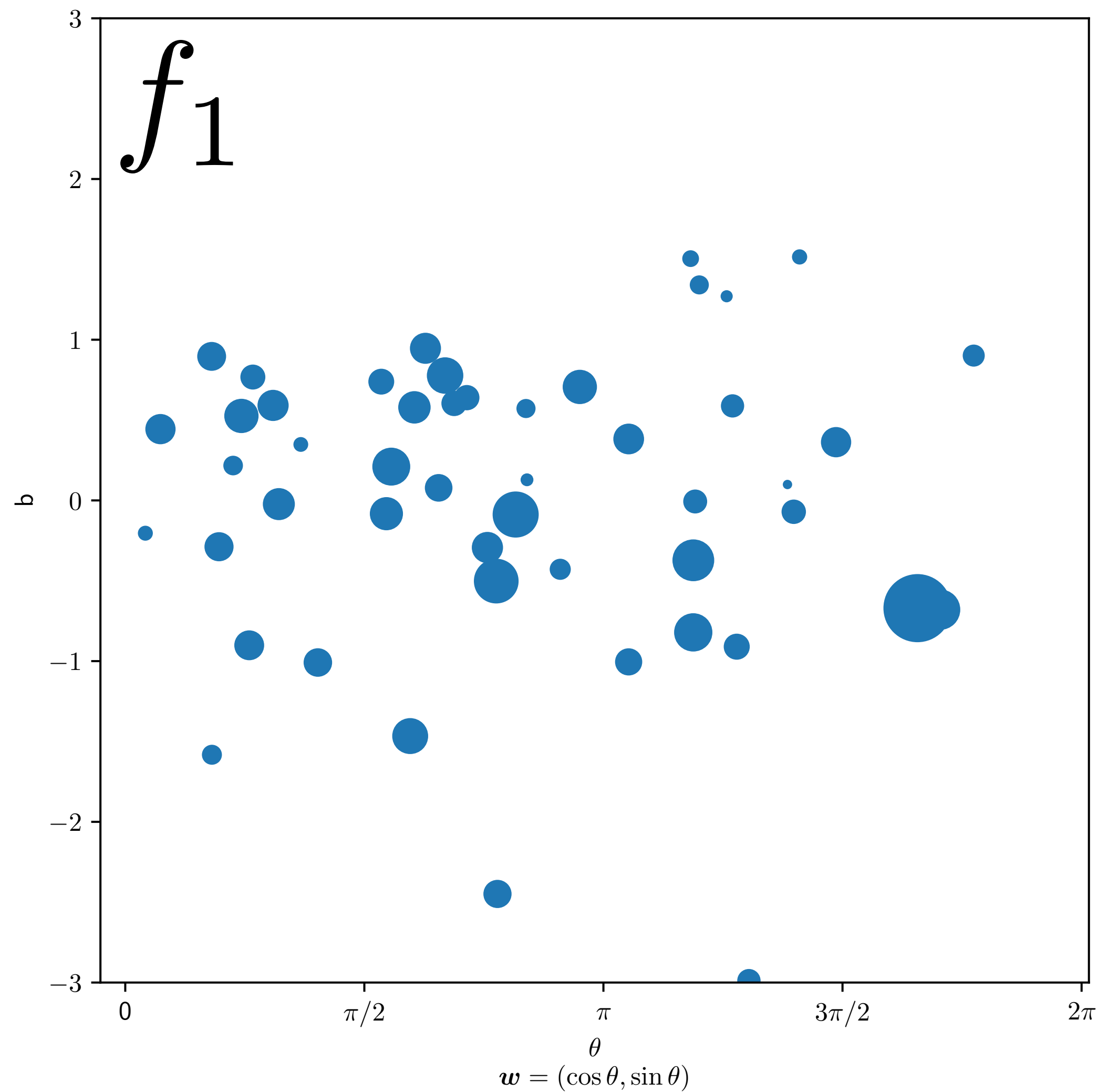
$$\|f\|_{VV} = O(\sqrt{D})$$

Simple Experiment

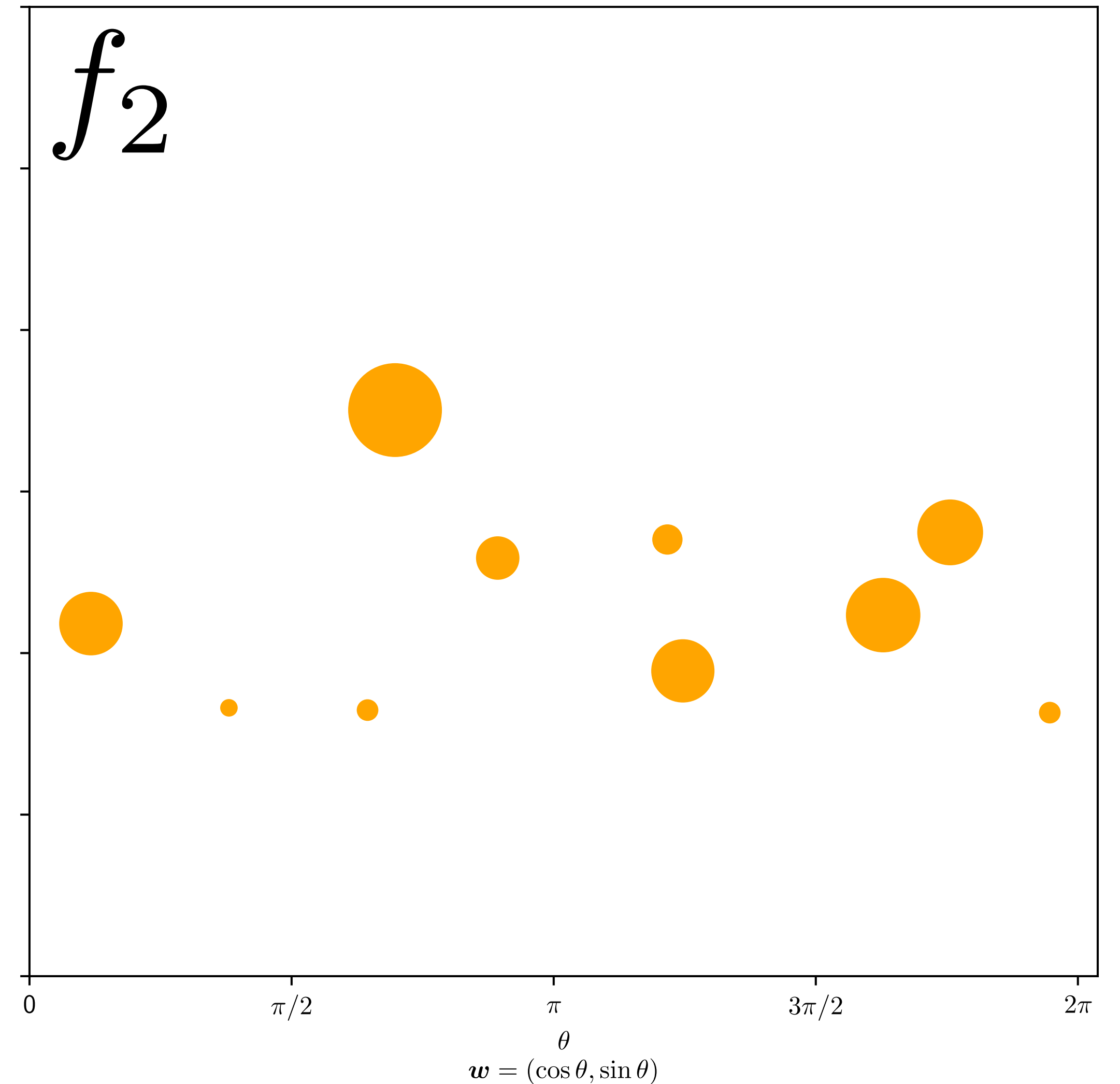
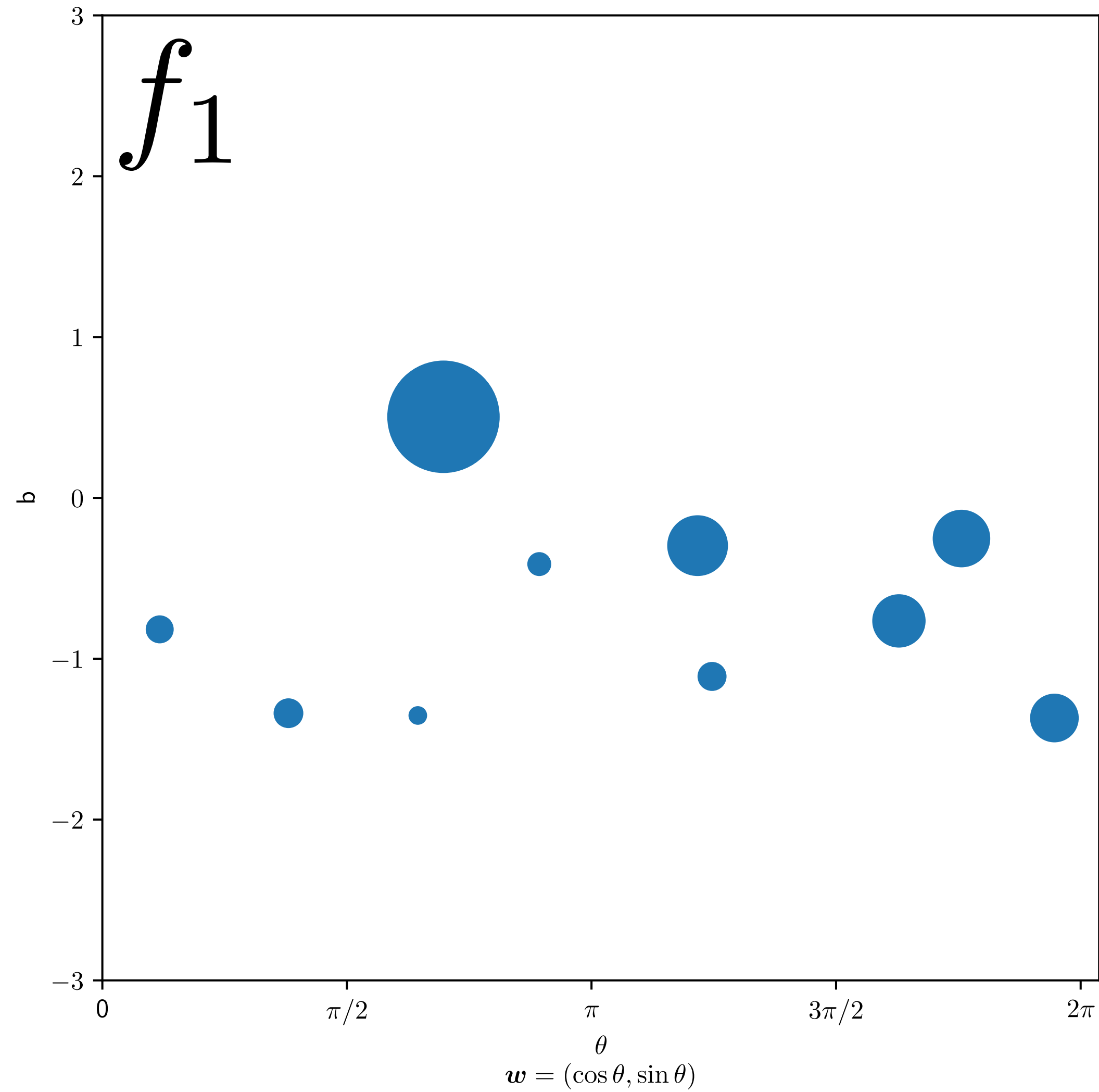


$$f(\mathbf{x}) = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \sum^K \mathbf{v}_k \sigma(\mathbf{w}_k^T \mathbf{x} - b_k)$$

No Weight Decay \rightarrow No Sharing



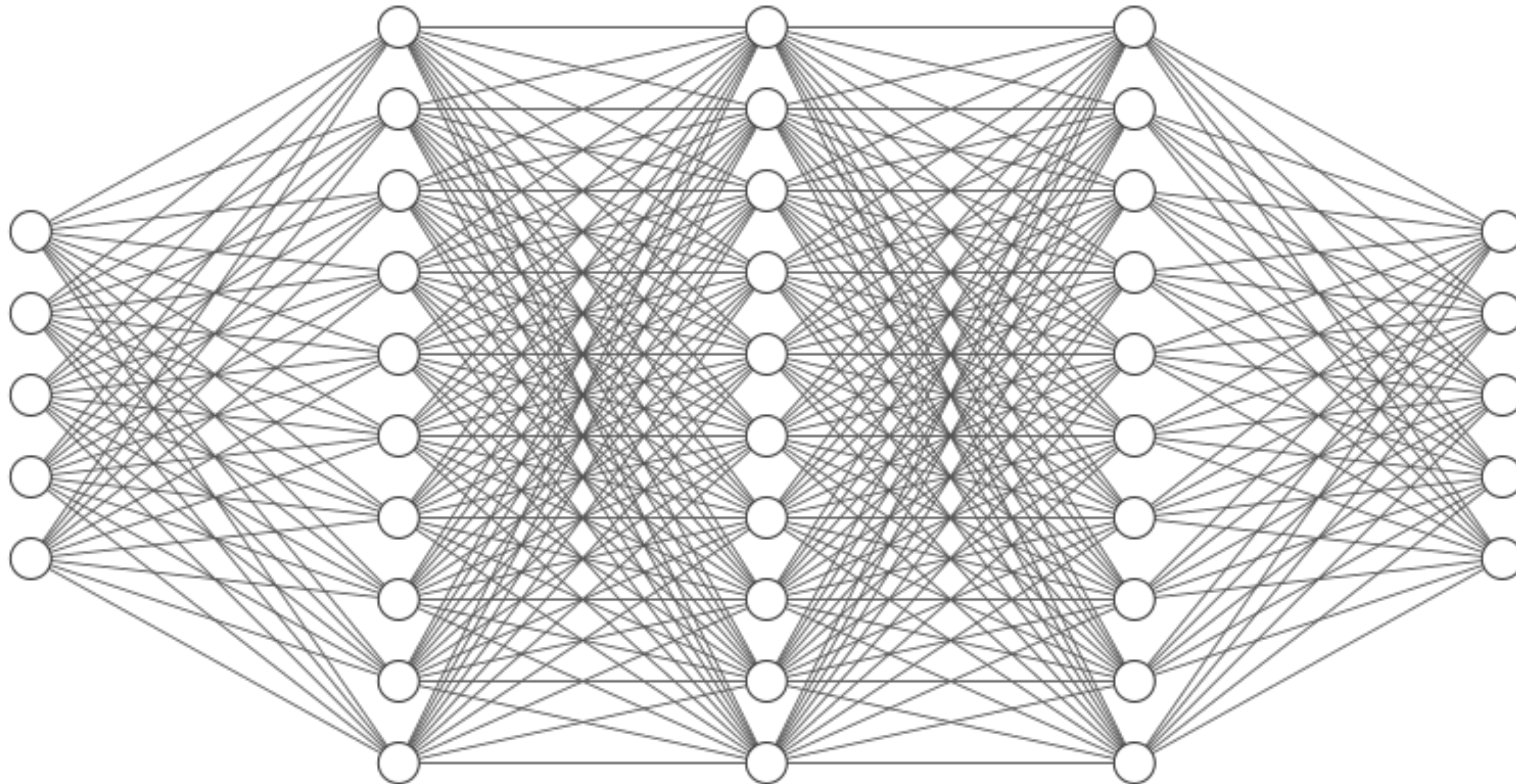
Weight Decay \rightarrow Neuron Sharing



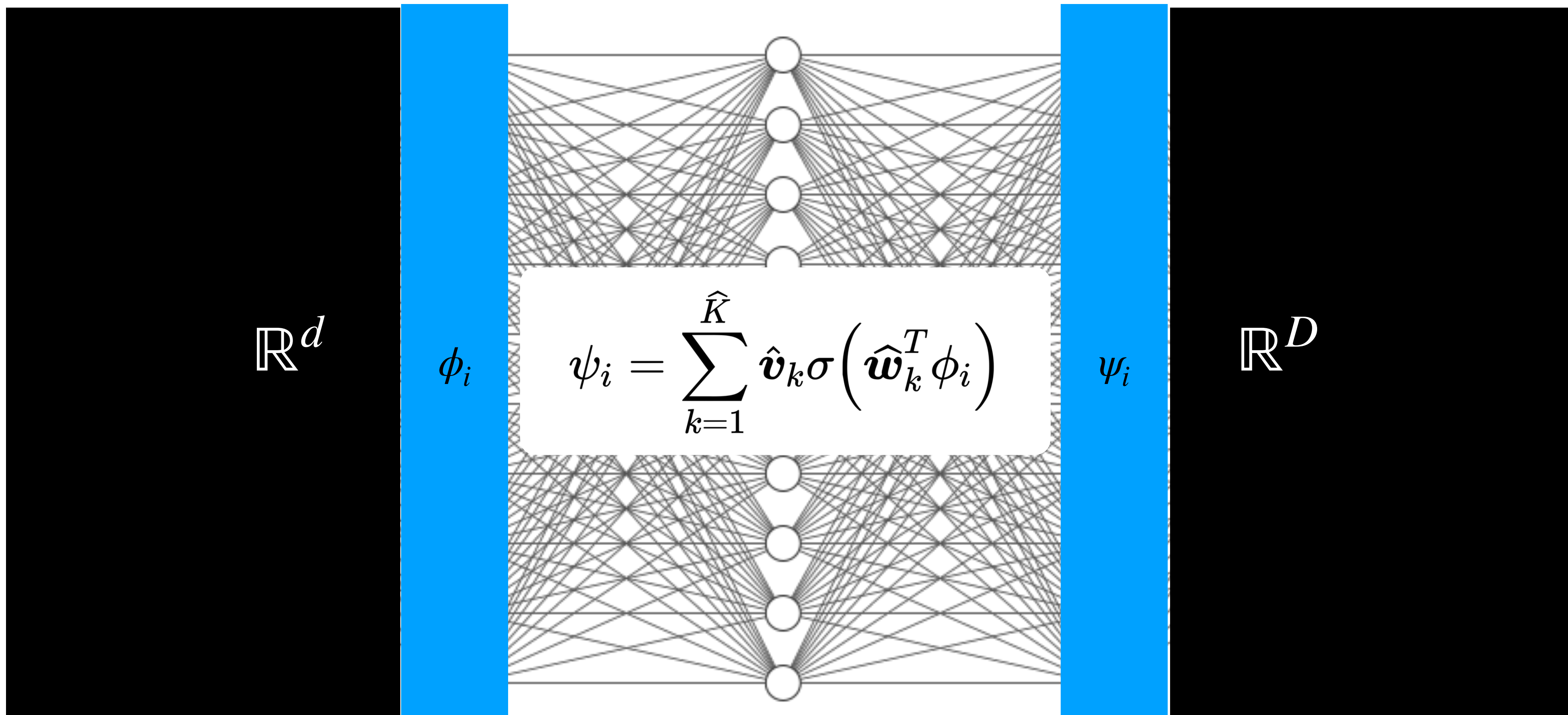
Tighter Bounds on Necessary Widths of DNNs

Bound on Network Width

- Suppose f_θ is a **DNN** which solves the weight decay objective



Bound on Network Width



$$\tilde{\Phi} = \left[\sigma(\hat{\mathbf{w}}_k^T \phi_1) \cdots \sigma(\hat{\mathbf{w}}_k^T \phi_N) \right]$$

$$\Psi = [\psi_1 \quad \psi_2 \cdots \psi_N]$$

Main Result

Assume f_θ is a DNN which solves the **weight decay objective**,

for any **homogenous layer** with input features ϕ_i and output features ψ_i , such that

$$\psi_i = \sum_{k=1}^{\hat{K}} \hat{\mathbf{v}}_k \sigma \left(\hat{\mathbf{w}}_k^T \phi_i \right) \quad i = 1, \dots, n$$

There exists another **optimal** representation of the form

$$\psi_i = \sum_{k=1}^K \mathbf{v}_k \sigma \left(\mathbf{w}_k^T \phi_i \right) \quad i = 1, \dots, n$$

Where,

$$K \leq \text{rank}(\Psi) \cdot \text{rank}(\tilde{\Phi})$$

Questions?

arXiv

